

# CUNY SPS DATA 621 - CTG5 - HW1

*Betsy Rosalen, Gabrielle Bartomeo, Jeremy O'Brien, Lidiia Tronina, Rose Koh*

*February 27, 2019*

## Contents

<b>1</b>	<b>DATA EXPLORATION</b>	<b>2</b>
1.1	Summary Statistics . . . . .	2
1.2	Shape of Predictor Distributions . . . . .	2
1.3	Outliers . . . . .	3
1.4	Missing Values . . . . .	3
1.5	Linearity . . . . .	3
<b>2</b>	<b>DATA PREPARATION</b>	<b>6</b>
2.1	Missing Values . . . . .	6
2.2	Remove Outliers . . . . .	6
2.3	Correlation . . . . .	6
2.4	Feature Engineering . . . . .	7
<b>3</b>	<b>BUILD MODELS</b>	<b>8</b>
3.1	MODEL 1 . . . . .	8
3.2	MODEL 2 . . . . .	10
3.3	MODEL 3 . . . . .	16
<b>4</b>	<b>SELECT MODELS</b>	<b>18</b>
4.1	Comparison of models . . . . .	18
4.2	Multi-collinearity . . . . .	18
4.3	Check Conditions for Least Squares Regression . . . . .	18
4.4	F-Statistic . . . . .	19
4.5	Predictions . . . . .	19
<b>5</b>	<b>Appendix</b>	<b>20</b>

---

# 1 DATA EXPLORATION

Professionals and gamblers alike are always seeking to optimize their chances of winning, whether it be sports, games, or their bets on them. Major League Baseball is a multibillion dollar industry where individual teams, players, and those who profit off of their success stand to benefit most from such optimization.

Data from 1871 to 2006 was collected in order to infer how many wins could be expected from the 162 games in a baseball team's season. Each observation represents a season for an unnamed team, and we have a total of 2,276 observations. For each team the target variable, TARGET\_WINS, represents the number of wins in a given year and has a maximum value of 162 possible wins. In addition to that 15 continuous integer predictor variables were collected (not including the index) representing each team's: base hits, doubles, triples, homeruns, walks, and strikeouts by batters, batters hit by pitches, bases stolen by batters and the number of times they were caught stealing, the number of errors, double plays, walks, hits, and homeruns allowed, and strikeouts by pitchers. The testing data contains the same 15 predictor variables and no target variable so it will be impossible to check the accuracy of our predictions from the testing data.

VARIABLE NAME	DEFINITION	THEORETICAL EFFECT ON WINS
TARGET_WINS	Number of wins	outcome variable
BATTING_H	Base Hits by batters (1B,2B,3B,HR)	Positive Impact
BATTING_2B	Doubles by batters (2B)	Positive Impact
BATTING_3B	Triples by batters (3B)	Positive Impact
BATTING_HR	Homeruns by batters (4B)	Positive Impact
BATTING_BB	Walks by batters	Positive Impact
BATTING_HBP	Batters hit by pitch (get a free base)	Positive Impact
BATTING_SO	Strikeouts by batters	Negative Impact
BASERUN_SB	Stolen bases	Positive Impact
BASERUN_CS	Caught stealing	Negative Impact
FIELDING_E	Errors	Negative Impact
FIELDING_DP	Double Plays	Positive Impact
PITCHING_BB	Walks allowed	Negative Impact
PITCHING_H	Hits allowed	Negative Impact
PITCHING_HR	Homeruns allowed	Negative Impact
PITCHING_SO	Strikeouts by pitchers	Positive Impact

## 1.1 Summary Statistics

Looking in Table 2, it can be easily noted that there are outliers present in more than one variable, with PITCHING\_H being the worst offender. Even at three times the standard deviation, its maximum value lays far outside of the 68-95-99.7 rule. FIELDING\_E, on the other hand, has the curious case of having a large difference between its mean and median, indicating there is skew present in this variable as well before any charts are actively looked at. Skewed variables cause bias in linear models and need treatment before being used.

## 1.2 Shape of Predictor Distributions

Figure. 1 shows the distribution of most of the variables seems normal although BASERUN\_SB, BASERUN\_CS, and BATTING\_3B have a slight to moderate right skew, FIELDING\_E, PITCHING\_BB, PITCHING\_H, and PITCHING\_SO have an extreme right skew, and BATTING\_HR, BATTING\_SO, and PITCHING\_HR are bimodal.

Table 2: Summary statistics

	n	min	mean	median	max	sd
TARGET_WINS	2276	0	80.79086	82.0	146	15.75215
BATTING_H	2276	891	1469.26977	1454.0	2554	144.59120
BATTING_2B	2276	69	241.24692	238.0	458	46.80141
BATTING_3B	2276	0	55.25000	47.0	223	27.93856
BATTING_HR	2276	0	99.61204	102.0	264	60.54687
BATTING_BB	2276	0	501.55888	512.0	878	122.67086
BATTING_SO	2174	0	735.60534	750.0	1399	248.52642
BASERUN_SB	2145	0	124.76177	101.0	697	87.79117
BASERUN_CS	1504	0	52.80386	49.0	201	22.95634
BATTING_HBP	191	29	59.35602	58.0	95	12.96712
PITCHING_H	2276	1137	1779.21046	1518.0	30132	1406.84293
PITCHING_HR	2276	0	105.69859	107.0	343	61.29875
PITCHING_BB	2276	0	553.00791	536.5	3645	166.35736
PITCHING_SO	2174	0	817.73045	813.5	19278	553.08503
FIELDING_E	2276	65	246.48067	159.0	1898	227.77097
FIELDING_DP	1990	52	146.38794	149.0	228	26.22639

As a result some data transformation will most likely be necessary to improve the accuracy of our model. The standard deviation of the various variables also hints at the intense skewing of some of the variables.

### 1.3 Outliers

Figure. 2 shows that there are also a large number of outliers that need to be accounted for, most prevalently in FIELDING\_E and BATTING\_H based off of the boxplots below. One such extreme outlier removed implied that there were, on average per game in a single season, 186 hits allowed by pitchers. This is an unrealistic figure, even for those for whom baseball is outside of their realm of understanding.

### 1.4 Missing Values

Figure. 3 displays of all the observations gathered across these fifteen variables, there are 3,478 missing values out of 36,416 total data points, which represents 10.187% of the data. Batters hit by pitches was missing the most, with 2,085 instances of missing information, which represents 91.61% of that variable missing. Additionally Pitching\_SO and Batting\_SO are missing exact same proportion 4.48% and are missing in the same observations. This data may not be missing at random and so there may be cause for removing it.

### 1.5 Linearity

Each variable was plotted against the target variable in order to determine at a glance which had the most potential linearity before the dataset was modified.

As can be observed in Figure. 4, the most influential variables are the ones previously discussed to have severe outliers and skew, and their linear relationship is negative - the higher the variable, the lower the target wins. On the other hand, BATTING\_H, BATTING\_BB and BATTING\_2B showed the most promise.

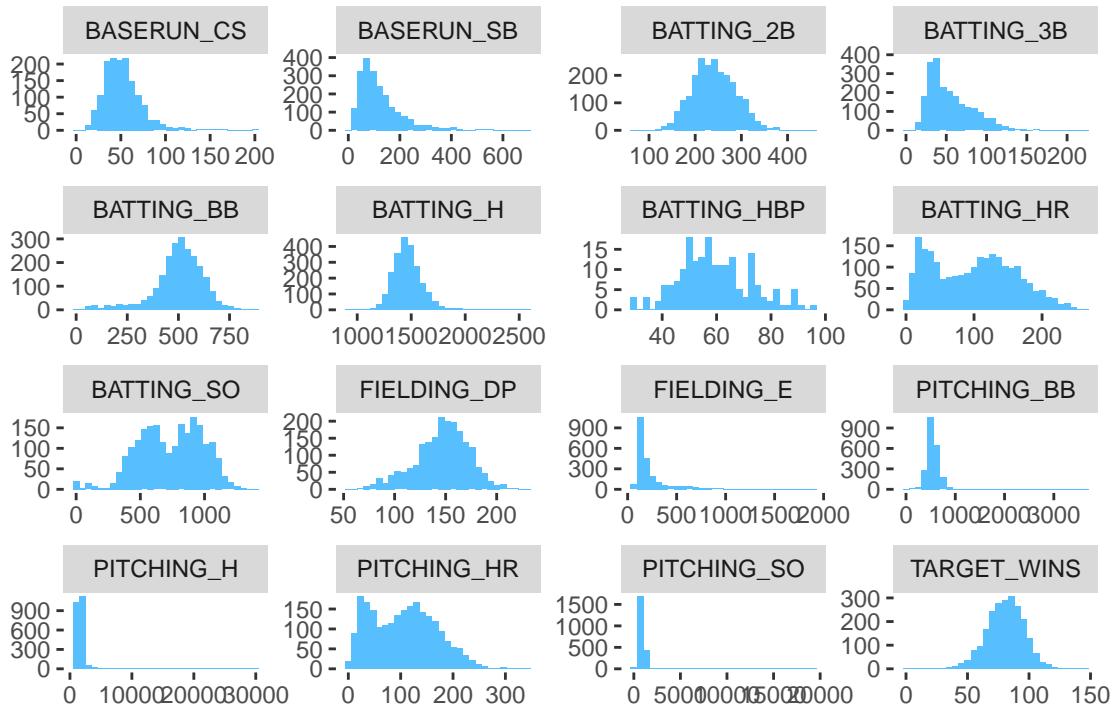


Figure 1: Data Distributions

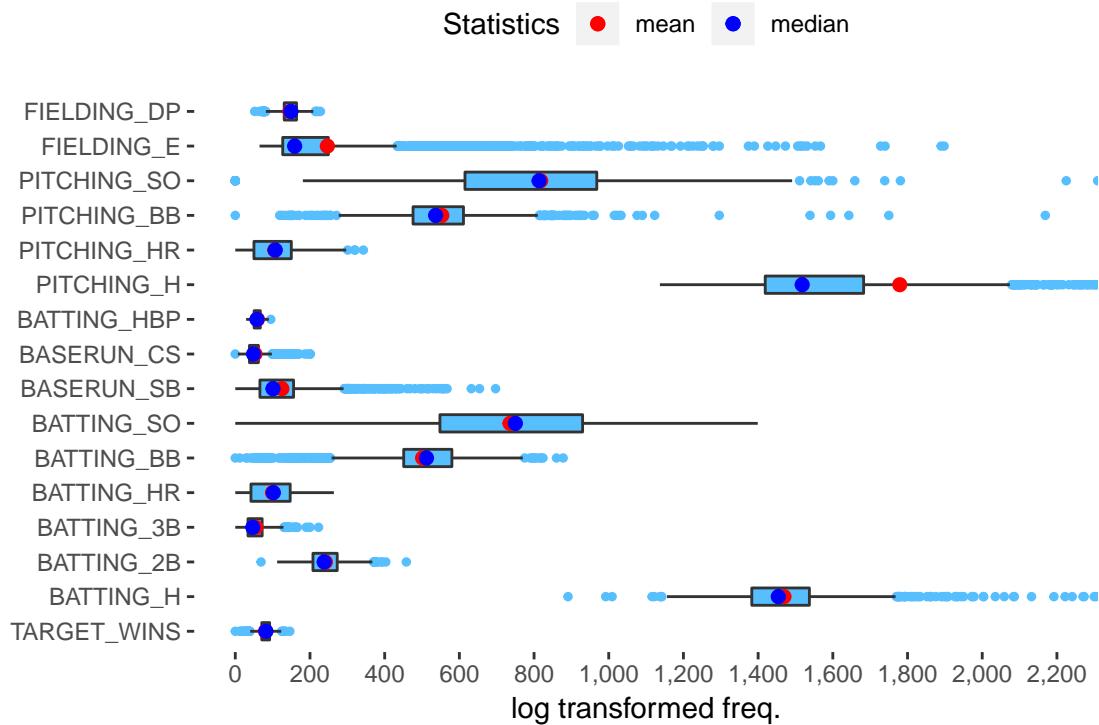


Figure 2: Boxplots highlighting many outliers in the data.

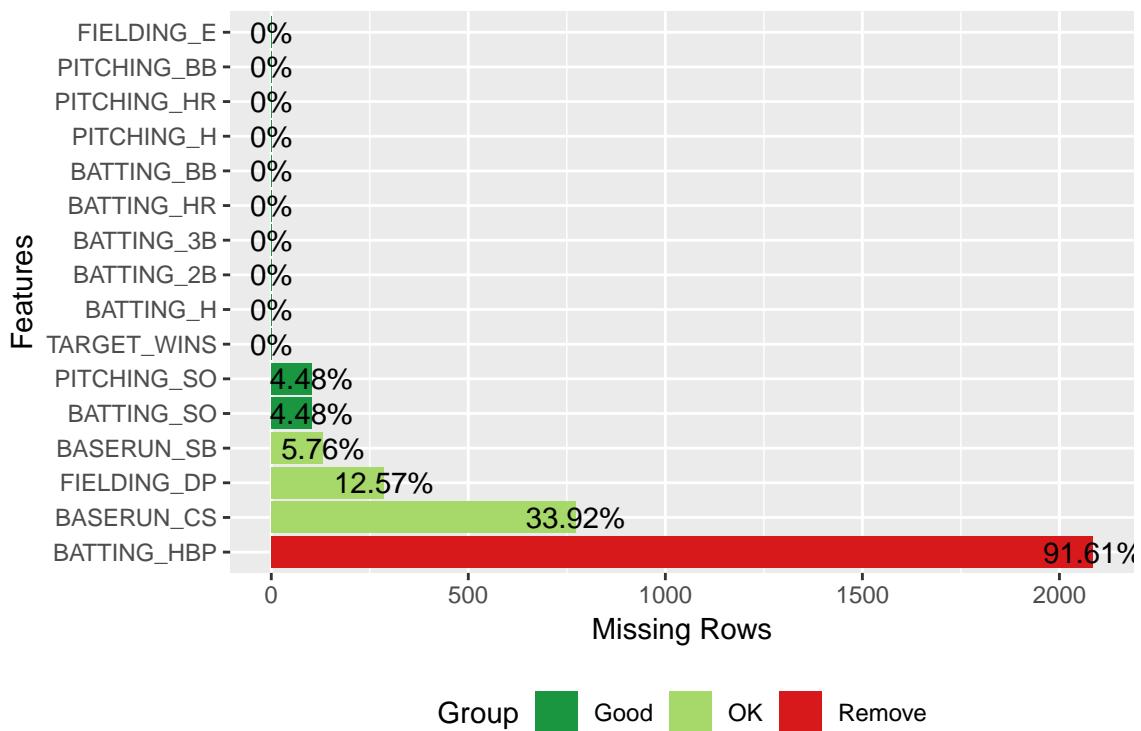


Figure 3: Missing values

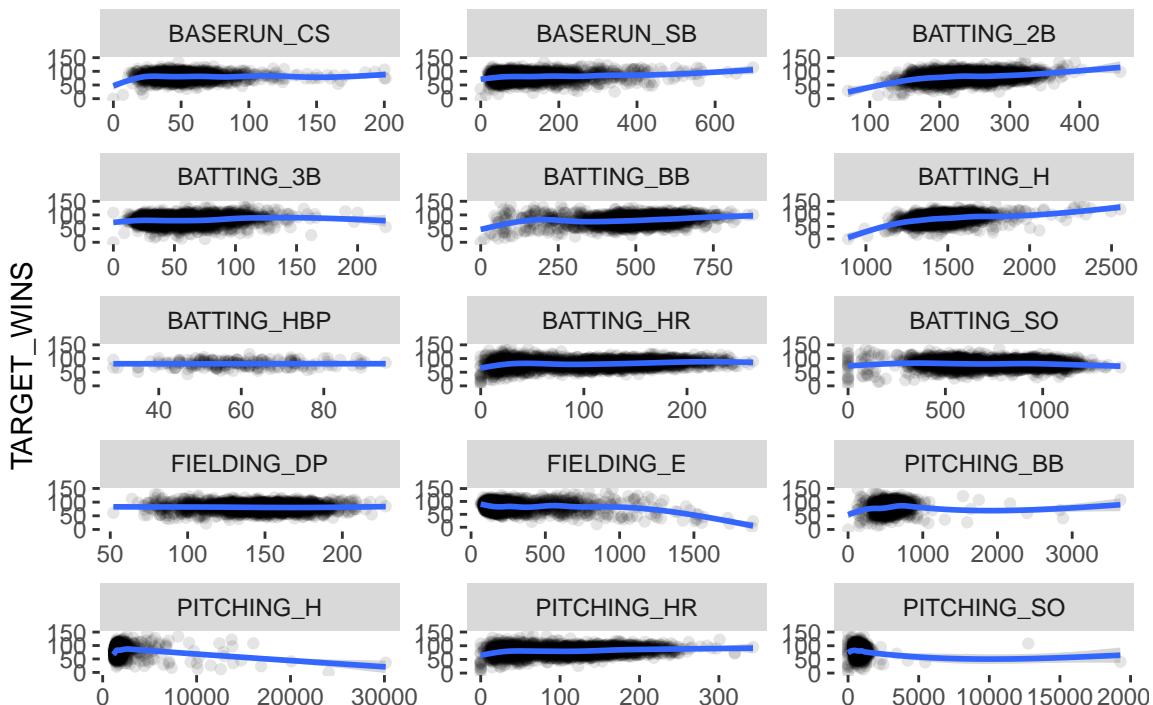


Figure 4: Linear relationships between each predictors and the target

---

## 2 DATA PREPARATION

### 2.1 Missing Values

As previously mentioned, just north of 10% of the data was missing values. Missing values can lead to errors in a model, bias, and worse if left unaccounted for. Attempting to “fix” this by imputing values or guessing why the values are missing in the first place - such as concluding that the missing values are meant to be zeroes - are just as likely to help with creating a model as it is to help with creating a disaster.

One of the R packages utilized, DataExplorer, which was used for Figure. 3, recommends removing null or missing values above a certain threshold as indicated in the graph.

Fixing missing values with imputation may help, but can also have a negative impact on the model if the assumed values do not correspond to the actual missing values. When it is just a few observations missing, modifications can be made, however, 91.61% is too large a proportion and would almost definitely distort the model, so we decided it was better to remove the `BATTING_HBP` column altogether. Deleting all cases with missing values, in this instance, would have shrunk the size of the dataset down to less than a tenth of its original size. If we simply delete all cases with missing values from the analysis, we will cause no bias, but we would most certainly lose a lot of important information.

Data that is Missing Completely at Random (MCAR), meaning the probability that a value is missing is the same for all cases can be imputed. Although there is some concern about whether or not `Pitching_S0` and `Batting_S0` are MCAR, we chose to leave all the remaining variables except `BATTING_HBP` and determine whether or not to remove them during the modelling process.

#### 2.1.1 NA Imputation

To deal with the remaining missing values, we used the `preProcess` function in the `caret` package. The `caret` package stands for Classification and regression training, it is a set of functions that streamline the process for creating predictive models, offers data splitting, pre-processing, feature selection, model tuning using resampling and variable importance estimationn.

For imputation, the selected method is `bagImpute`. For each predictor in the data, a bagged tree is created using all of the other predictors in the train dataset. When a new sample has a missing predictor value, the bagged model is used to predict the value. In theory, this is a more powerful method of imputing compared to `KnnImpute`, however, much higher computational costs are a downside. We can see the difference between the original and imputed data in Table 3.

### 2.2 Remove Outliers

Outlier treament was done by placing a threshold of five times the standard deviation up from the mean and removing all observations that fell north of this boundary.

### 2.3 Correlation

The theoretical effect of strikeouts by batters, batters caught stealing, errors, walks, hits, and homeruns allowed were believed to have a negative impact on the number of wins of an individual team in a given year. A closer look at the correlation plot between the variables painted a different picture. Figure. 5 shows correlation plot.

Table 3: Difference between original and imputed data

	n	min	mean	median	max	sd
TARGET_WINS	0	0	0.0000000	0	0	0.000000
BATTING_H	0	0	0.0000000	0	0	0.000000
BATTING_2B	0	0	0.0000000	0	0	0.000000
BATTING_3B	0	0	0.0000000	0	0	0.000000
BATTING_HR	0	0	0.0000000	0	0	0.000000
BATTING_BB	0	0	0.0000000	0	0	0.000000
BATTING_SO	-102	0	7.1015741	11	0	2.031465
BASERUN_SB	-131	0	-0.6324276	-4	0	2.342842
BASERUN_CS	-772	0	-19.2926570	-8	0	-17.471877
PITCHING_H	0	0	0.0000000	0	0	0.000000
PITCHING_HR	0	0	0.0000000	0	0	0.000000
PITCHING_BB	0	0	0.0000000	0	0	0.000000
PITCHING_SO	-102	0	10.0065202	10	0	10.012499
FIELDING_E	0	0	0.0000000	0	0	0.000000
FIELDING_DP	-286	0	0.9773713	2	0	1.306904

When compared to what was hypothesized, there was actually a positive impact for the number of wins for a team in a given year by walks, hits, and homeruns allowed; at the same time, variables previously thought to have a positive correlation - strikeouts by pitchers and double plays - had a negative correlation for the number of wins. The three variables with the greatest correlation to the number of wins were the hits allowed, the walks by batters, and the walks allowed. Of these, the hits allowed had a relatively low correlation with the walks by batters and the walks allowed, whereas the walks allowed and the walks by batters had a direct positive correlation with one another.

---

## 2.4 Feature Engineering

Since there are four pairs of related variables that are two sides of the same coin, hits allowed vs. hits by batters, home runs allowed vs. home runs hit by batters, etc. and three of those pairs are highly correlated with each other we decided to try using the difference between them in place of the original variables in our original models. We decided to use offense (batting) minus defense (pitching). These arithmetically transformed offense / defense variables are linearly related with BATTING and PITCHING variables, so we can include one or the other in a model, but not both. However, replacing original variables with these transforms did not improve  $R^2$  in a base case.

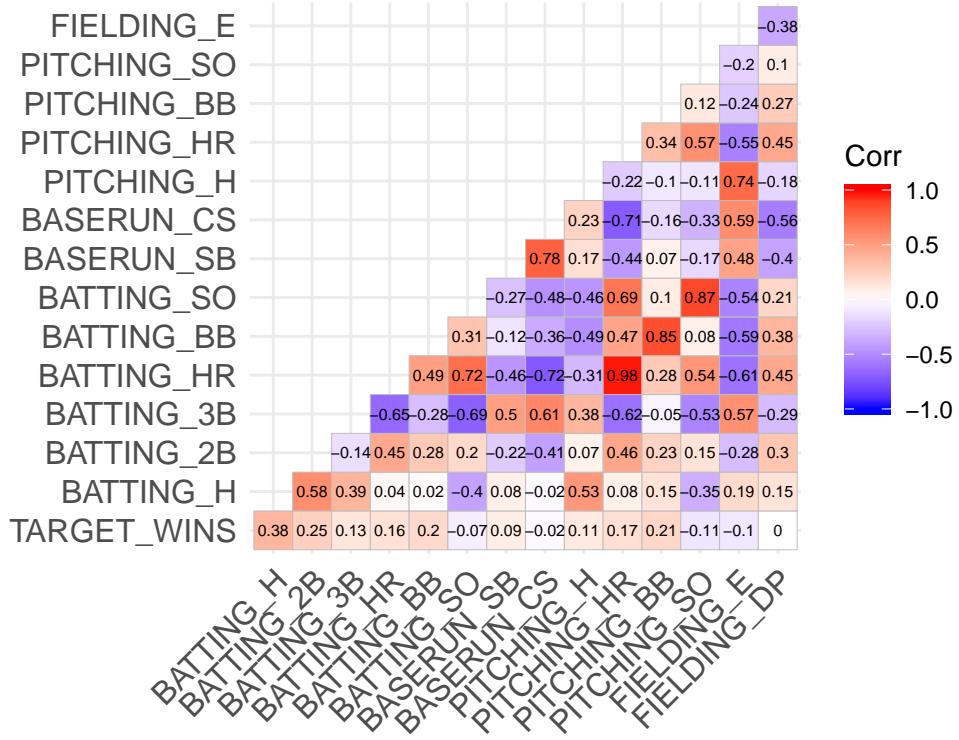


Figure 5: Correlation

### 3 BUILD MODELS

#### 3.1 MODEL 1

Multiple regression can be created as a purely statistical model through the use of significance tests, or it can be interpreted in a more practical, non-statistical manner. This first approach is more of the latter and is based on consultation with a subject-area expert.

After discussing the various valuables available with the expert, the following categories from the most important to the least important variables were developed:

**Very Important:** BATTING\_H, BATTING\_HR, BATTING\_SO, FIELDING\_E, PITCHING\_SO

**Fairly Important:** BASERUN\_SB, PITCHING\_HR, BATTING\_BB

**Important:** BATTING\_2B, BATTING\_3B, FIELDING\_DP, PITCHING\_H

**Slightly Important:** PITCHING\_BB, BASERUN\_CS

**Not Important:** BATTING\_HBP

Reviewing this preliminary categorization led to the expert suggesting the elimination of BATTING\_HBP and BASERUN\_CS as they are the least important when determining the number of wins in a season as per their professional opinion.

```

##             Estimate Std. Error     t value   Pr(>|t|) 
## (Intercept) 36.96752362 5.659839939  6.5315493 8.047746e-11
## BATTING_H    0.01636094 0.005009596  3.2659206 1.107724e-03
## BATTING_HR   0.10817201 0.050561725  2.1394052 3.251183e-02
## BATTING_SO   -0.02840724 0.006545850 -4.3397329 1.490644e-05
## FIELDING_E   -0.03883797 0.003444021 -11.2769272 1.017041e-28
## PITCHING_SO  0.01723395 0.005381052  3.2027097 1.380717e-03
## BASERUN_SB   0.04307539 0.004943029  8.7143716 5.616586e-18
## PITCHING_HR  -0.02312490 0.046842474 -0.4936737 6.215855e-01
## BATTING_BB   0.07700679 0.016212347  4.7498854 2.165548e-06
## BATTING_2B   -0.01136604 0.009372708 -1.2126740 2.253837e-01
## BATTING_3B   0.12831020 0.018090002  7.0928796 1.759636e-12
## FIELDING_DP  -0.09303898 0.013066848 -7.1202314 1.449816e-12
## PITCHING_BB  -0.05442930 0.014378645 -3.7854263 1.575214e-04
## PITCHING_H   0.01368874 0.001538180  8.8993081 1.141474e-18

```

Creating a linear model sans the two aforementioned variables led to an Adjusted  $R^2$  of 0.2793 on Adjusted  $R^2$ . The  $R^2$ 's value decreased when attempting to remove other less-than-important variables.

While forward selection had mild success, it was determined the best next step would be to perform backwards elimination. Based on this method and its results, BATTING\_H and BATTING\_2B were removed.

This resulted in the following model:

```

TARGET_WINS ~ BATTING_H + BATTING_HR + BATTING_SO +
  FIELDING_E + PITCHING_SO + BASERUN_SB + BATTING_BB + BATTING_3B +
  FIELDING_DP + PITCHING_BB + PITCHING_H

##             Estimate Std. Error     t value   Pr(>|t|) 
## (Intercept) 51.02945220 3.804274044 13.4137162 1.667309e-39
## BATTING_HR   0.12517282 0.050392683  2.4839483 1.306679e-02
## BATTING_SO   -0.03037268 0.006527324 -4.6531593 3.461565e-06
## FIELDING_E   -0.04213970 0.003250397 -12.9644760 4.184217e-37
## PITCHING_SO  0.01592382 0.005340224  2.9818633 2.896198e-03
## BASERUN_SB   0.04739597 0.004787681  9.8995676 1.222747e-22
## PITCHING_HR  -0.02606546 0.046927513 -0.5554409 5.786491e-01
## BATTING_BB   0.08441157 0.016037786  5.2632934 1.551180e-07
## BATTING_3B   0.15744729 0.015968741  9.8597185 1.789587e-22
## FIELDING_DP  -0.08569435 0.012908313 -6.6386947 3.966504e-11
## PITCHING_BB  -0.06060056 0.014257515 -4.2504297 2.221945e-05
## PITCHING_H   0.01700111 0.001182808  14.3735124 7.370476e-45

```

The  $R^2$  produced was still low (0.276), so outliers were analyzed as more thoroughly due to the effect they can have on models. PITCHING\_H had a high number of outliers which indicated a need for data transformation; log transformation was used to normalize this variable and further bring the linear model in line.

```

##             Estimate Std. Error     t value   Pr(>|t|) 
## (Intercept) -3.043824e+02 24.216098678 -12.5694254 4.724977e-35
## BATTING_HR    1.020673e-01 0.049941697  2.0437299 4.109815e-02
## BATTING_SO   -1.245250e-02 0.006620296 -1.8809574 6.010862e-02
## FIELDING_E   -4.236317e-02 0.003162634 -13.3949008 2.108130e-39
## PITCHING_SO   4.800476e-03 0.005296872  0.9062851 3.648834e-01
## BASERUN_SB    4.275539e-02 0.004682297  9.1312870 1.480875e-19
## PITCHING_HR  -3.182162e-02 0.046363949 -0.6863440 4.925679e-01

```

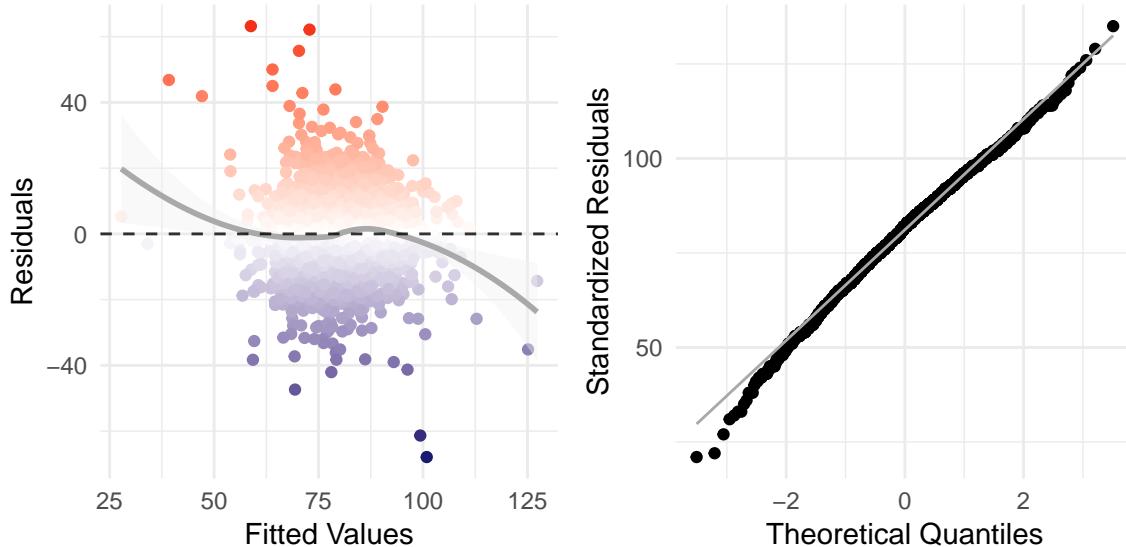


Figure 6: Model 1: Residual Plot and Q-Q Plot

```
## BATTING_BB      9.903885e-02  0.015976185   6.1991552 6.749227e-10
## BATTING_3B     1.180135e-01  0.015866324   7.4379893 1.452989e-13
## FIELDING_DP    -9.572525e-02  0.012786300  -7.4865479 1.014106e-13
## PITCHING_BB    -7.367309e-02  0.014212070  -5.1838394 2.370449e-07
## log(PITCHING_H) 5.231226e+01  3.236934315  16.1610501 1.258201e-55
```

After we used the log transformation the model's Adjusted  $R^2$  increased to 0.2961.

Residuals for this linear model are normally distributed and random. The Q-Q plot confirms this model can be used to predict the number of wins in a season for a team, though there is still room for improvement using other methods.

Below is Model 1's prediction result for the test data:

```
##      Min. 1st Qu. Median Mean 3rd Qu. Max. NA's
## 60.74 75.98 81.12 81.43 86.03 108.43 54
```

### 3.1.1 Summary of Results

r.squared	adj.r.squared	statistic	p.value	df	deviance	df.residual
0.2976822	0.2941991	85.46496	0	12	346538.5	2218

It was determined the overall subject-area expertise wasn't as effective as a stand-alone method of creating multiple regression models. Statistical iterations which were performed contradicted the subject area expert such as removing `BATTING_H` from the model. Additionally the log transformation of `PITCHING_H` made a significant improvement in the model's linearity.

## 3.2 MODEL 2

Our approach for Model 2 was to try to use as many of the tools as possible that are available in R and that we have learned thus far to determine a model. This was based solely on the statistical qualities of the

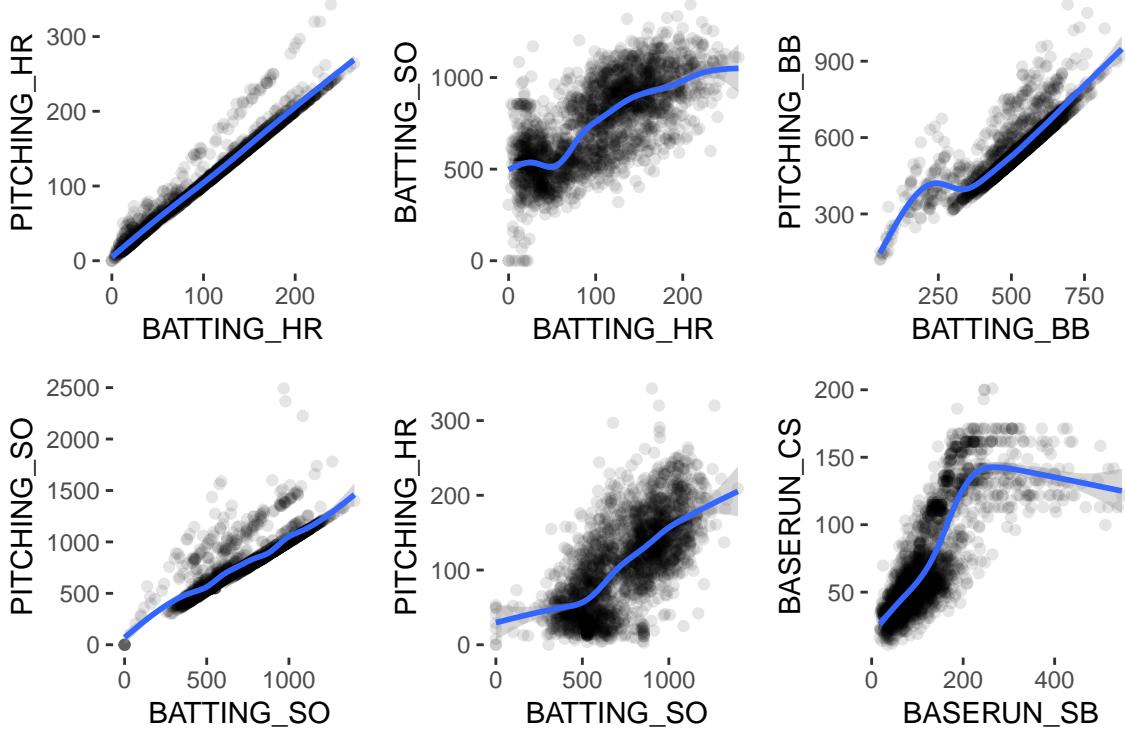


Figure 7: Scatterplots showing possible collinearity problems

predictor variables without any regard to our expert's opinion.

### 3.2.1 Check for Correlated Predictor Variables and Linear Relationship to Target

We started by plotting the relationships between variables that had high correlation values to look for potential collinearity problems.

Based on the Figure. 6, we decided somewhat arbitrarily to remove the three pitching variables (PITCHING\_HR, PITCHING\_BB, and PITCHING\_SO) rather than the corresponding batting variables (BATTING\_HR, BATTING\_BB, and BATTING\_SO) due to the extremely high correlation between these predictors.

We then plotted the remaining variables to see if they showed a linear relationship with the target variable. Most of the remaining predictors showed a clear linear relationship with the target, however, the extreme skew of PITCHING\_H and FIELDING\_E as well as a more moderate skew in BASERUN\_SB and BATTING\_3B, can be seen in the plots.

### 3.2.2 Log Transform Data

We decided to log transform PITCHING\_H, FIELDING\_E, BASERUN\_SB and BATTING\_3B in order to compensate for the skew. The resulting distributions can be seen in the revised plots in Figure. 7.

### 3.2.3 Building the First Model

Finally we built a model based on the selected variables including the log transformations where appropriate.

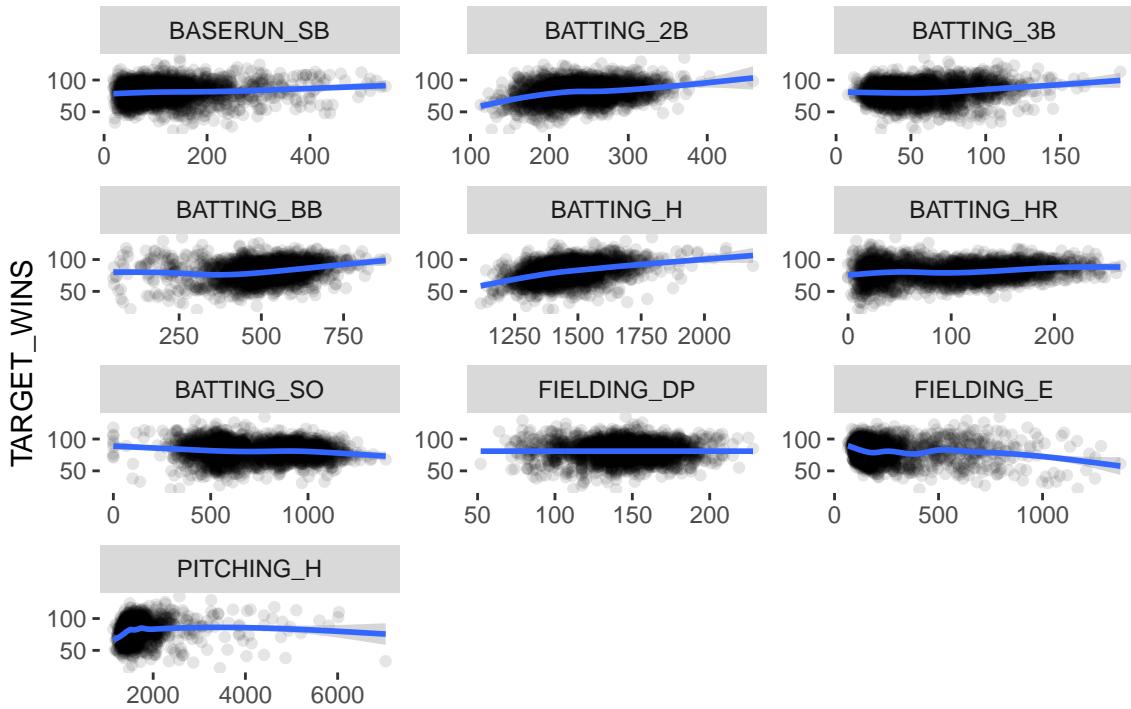


Figure 8: Linear relationship between each predictor and the target showing highly skewed variables

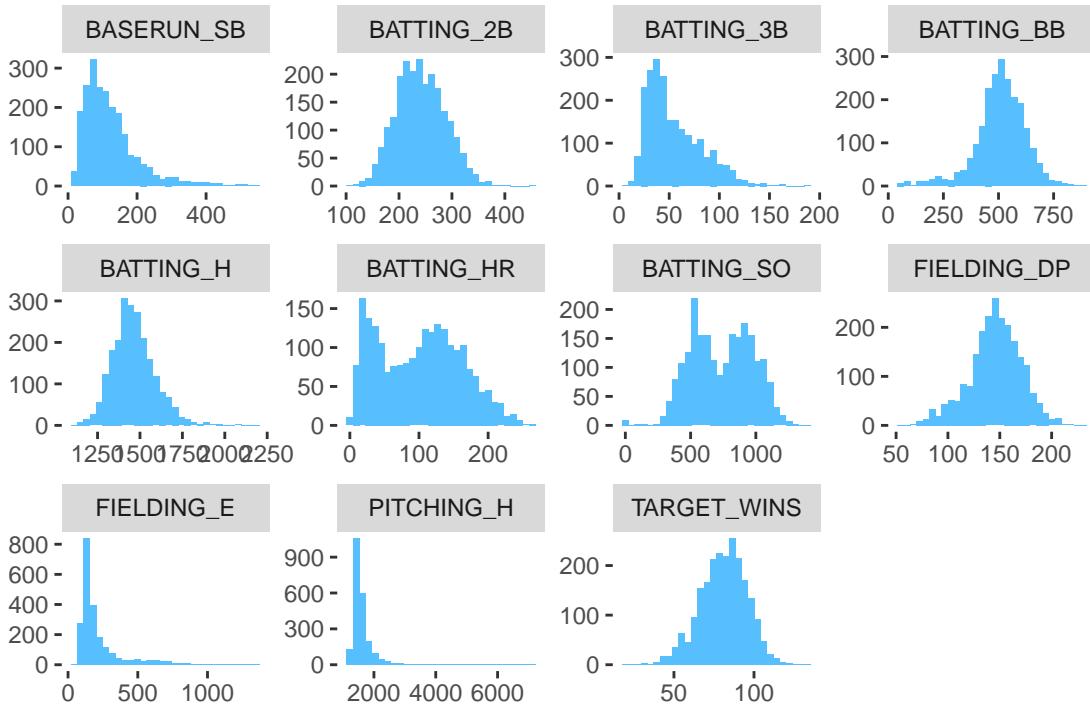


Figure 9: Predictor variable distributions showing highly skewed variables

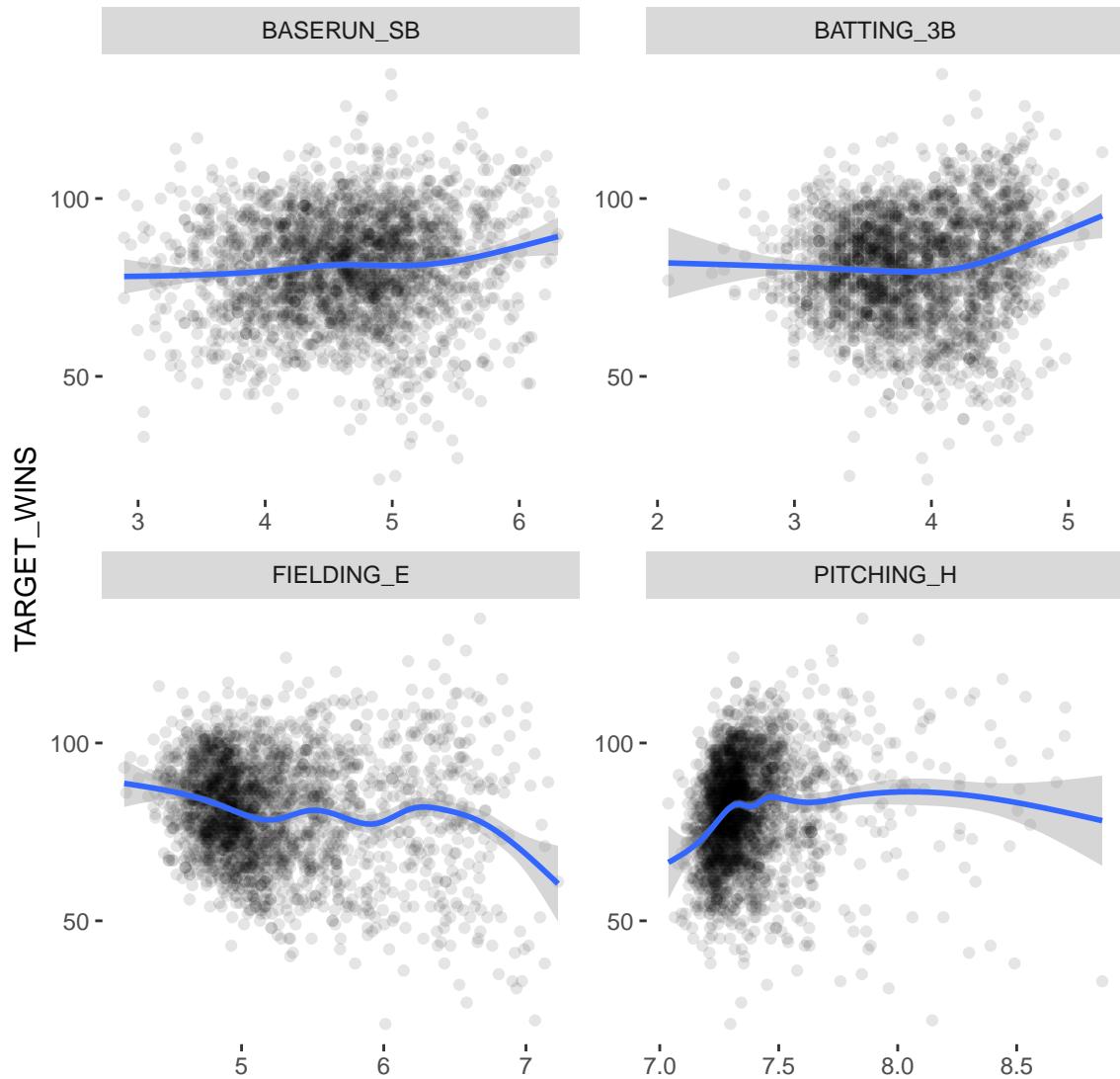


Figure 10: Linear relationship between each log transformed predictor and the Target showing decreased skew

### Histograms

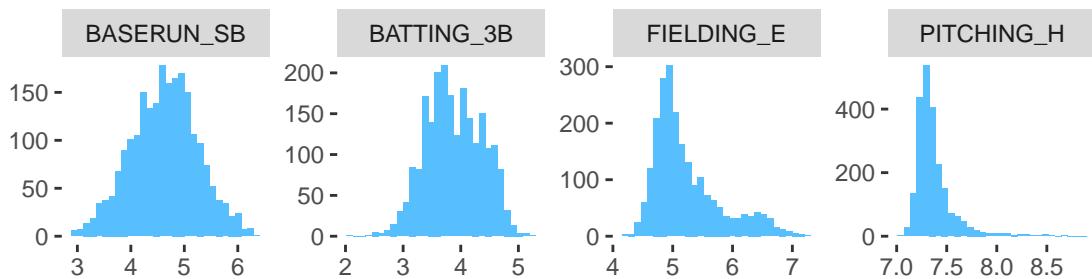


Figure 11: Log transformed distributions showing decreased skew

Table 4: First Model Coefficients

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-60.4527921	15.8795234	-3.806965	0.0001445
BATTING_H	0.0261430	0.0045993	5.684135	0.0000000
BATTING_2B	-0.0310908	0.0091945	-3.381443	0.0007336
log(BATTING_3B)	7.0262602	0.9457691	7.429150	0.0000000
BATTING_HR	0.0718230	0.0101324	7.088438	0.0000000
BATTING_BB	0.0195462	0.0031847	6.137572	0.0000000
BATTING_SO	-0.0115911	0.0024155	-4.798646	0.0000017
log(BASERUN_SB)	4.7452615	0.5628468	8.430822	0.0000000
log(PITCHING_H)	18.6870643	2.6087465	7.163235	0.0000000
log(FIELDING_E)	-13.2410442	1.0626760	-12.460095	0.0000000
FIELDING_DP	-0.1079247	0.0131023	-8.237083	0.0000000

```
TARGET_WINS ~ BATTING_H + BATTING_2B + log(BATTING_3B) + BATTING_HR +
  BATTING_BB + BATTING_SO + log(BASERUN_SB) + log(PITCHING_H) +
  log(FIELDING_E) + FIELDING_DP
```

All of the variables had a very low p-value indicating a significant impact on our target, however our  $R^2$  value was low at only 0.2889.

### 3.2.3.1 First Model $R^2$ 0.2909412

### 3.2.4 Refining the Model with leaps Package

We thought we may be able to use some other tools in R to refine our model and get a better  $R^2$  value. So next we tried using the leaps package to see if it would recommend removing any of our chosen variables from the model. In the following plot you can see that we could remove BATTING\_H, BATTING\_2B without affecting out  $R^2$  much, but it would not improve the model.

```
## NULL
```

### 3.2.5 Refining the Model by Standardizing the Predictor Variables

Next we tried standardizing the (non-log-transformed) variables to see what impact that might have on our model. Standardizing actually resulted in a significant reduction in our  $R^2$  value from 0.2889 to 0.274.

### 3.2.5.1 STANDARDIZED Model $R^2$ 0.2768969

### 3.2.6 Test all of the predictors

Next we ran an ANOVA test to compare our model to the null model. With a p-value that is basically zero, clearly our model is statistically significant.

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
2229	493421.2	NA	NA	NA	NA
2219	349864.7	10	143556.5	91.05005	0

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
2220	360087.7	NA	NA	NA	NA
2219	349864.7	1	10223	64.83889	0

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
2220	357027.1	NA	NA	NA	NA
2219	349864.7	1	7162.371	45.427	0

### 3.2.7 Testing a subspace

We then tried testing a subspace. Since our initial models using the difference between the corresponding batting and pitching variables did not show promise we tried adding those two variables instead.

Once again our model declined in performance rather than improving.

#### 3.2.7.1 Subspace Model $R^2$ 0.289032

### 3.2.8 Refining the Model with the MASS Package

Last, but not least, we used the stepAIC function from the MASS package to see if it came up with different recommendations for what variabels to keep and which to exclude from our model. We started with all variables putting back the ones we had previously taken out due to collinearity issues and let the algorithm choose which to keep.

The final suggested model was:

```
Final Model:
TARGET_WINS ~ BATTING_H + BATTING_3B + BATTING_HR + BATTING_BB +
    BATTING_SO + BASERUN_SB + BASERUN_CS + PITCHING_H + PITCHING_BB +
    PITCHING_SO + FIELDING_E + FIELDING_DP
```

In comparison to our original model we had the following variables added to our model (BASERUN\_CS, PITCHING\_BB, and PITCHING\_SO) and the following variable removed (BATTING\_2B).

We tried multiple iterations of that model, without any log transformations, with log transformations, with and without the collinear variables, but whenever we removed one of the collinear variables our model would decline in performance, so we decided to try our multiplying the corresponding collinear variables together and BINGO! We got an  $R^2$  of 0.3247 using the following model:

```
TARGET_WINS ~ BATTING_3B + BATTING_HR + BATTING_BB*PITCHING_BB +
    BATTING_SO*PITCHING_SO + BASERUN_SB + BASERUN_CS + BATTING_H*log(PITCHING_H) +
    log(FIELDING_E) + FIELDING_DP
```

#### 3.2.8.1 FINAL Model 2 $R^2$ 0.3273767

---

Table 5: FINAL Model 2 Coefficients

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	184.7024691	103.5096997	1.7843977	0.0744960
BATTING_3B	0.1609100	0.0178183	9.0306083	0.0000000
BATTING_HR	0.0726106	0.0101624	7.1450265	0.0000000
BATTING_BB	0.0579817	0.0198914	2.9149165	0.0035936
PITCHING_BB	-0.0827055	0.0129714	-6.3759862	0.0000000
BATTING_SO	0.0072940	0.0091561	0.7966181	0.4257583
PITCHING_SO	0.0252878	0.0064281	3.9339378	0.0000861
BASERUN_SB	0.0302723	0.0060897	4.9710708	0.0000007
BASERUN_CS	0.0619866	0.0151277	4.0975680	0.0000433
BATTING_H	-0.2119680	0.0542250	-3.9090439	0.0000954
log(PITCHING_H)	-5.4137306	14.2877118	-0.3789082	0.7047924
log(FIELDING_E)	-16.6786573	1.1073889	-15.0612470	0.0000000
FIELDING_DP	-0.0901744	0.0131173	-6.8744807	0.0000000
BATTING_BB:PITCHING_BB	0.0000473	0.0000168	2.8074726	0.0050369
BATTING_SO:PITCHING_SO	-0.0000284	0.0000044	-6.4798217	0.0000000
BATTING_H:log(PITCHING_H)	0.0293581	0.0074025	3.9659810	0.0000754

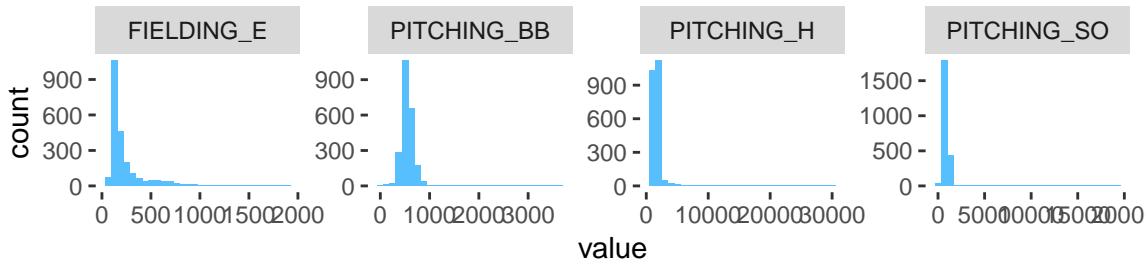


Figure 12: Histograms of variables showing pronounced rightward-skew

### 3.3 MODEL 3

We sought to explore whether there was a relationship between wins and the difference of specific offensive and defensive team capabilities - hits, homeruns, balls, and strike-outs. Incorporating variables that reflect those differences (i.e. subtracting batting hits from pitching hits, and so on), however, did not improve the explanatory power of the model beyond using the original variables.

Given these variables did not yield improvements, in their place we explored a third model. As the histograms below highlight, a number of the independent variables - pitching hits, pitching homeruns, pitching strikeouts - demonstrate pronounced rightward-skew.

We corrected that skew by transforming those three variables using natural logarithms. When we tested those log transformations in a model where they replaced the untransformed original variables combined with all other variables, we found that neither the originals nor the log transformations for pitching homeruns and pitching strikeouts met the threshold of significance (a p-value below the  $\alpha$  level of .05). Based on high p-values, over a series of backward steps we removed pitching homeruns, pitching strikeouts, and baserun caught stealing, yielding the following model:

#### 3.3.1 Residual and Q-Q Plot

Based on this model's F-statistic and p-value, we can reject the null hypothesis that coefficients with values of zero would fit the data better. Per the adjusted  $R^2$  value, this model explains approximately 29.56% of

Table 6: Log Transform Model Coefficients

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-405.5737568	36.8799143	-10.997145	0.0000000
BATTING_H	-0.0134073	0.0060215	-2.226592	0.0260746
BATTING_2B	-0.0159922	0.0091168	-1.754150	0.0795429
BATTING_3B	0.1411944	0.0175995	8.022649	0.0000000
BATTING_HR	0.0741949	0.0097120	7.639496	0.0000000
BATTING_BB	0.1254287	0.0126385	9.924335	0.0000000
BATTING_SO	-0.0070022	0.0023309	-3.004051	0.0026939
BASERUN_SB	0.0460848	0.0047847	9.631699	0.0000000
log(PITCHING_H)	68.7851289	5.7418499	11.979611	0.0000000
PITCHING_BB	-0.0958116	0.0106148	-9.026255	0.0000000
FIELDING_E	-0.0484349	0.0035654	-13.584555	0.0000000
FIELDING_DP	-0.0897449	0.0128950	-6.959664	0.0000000

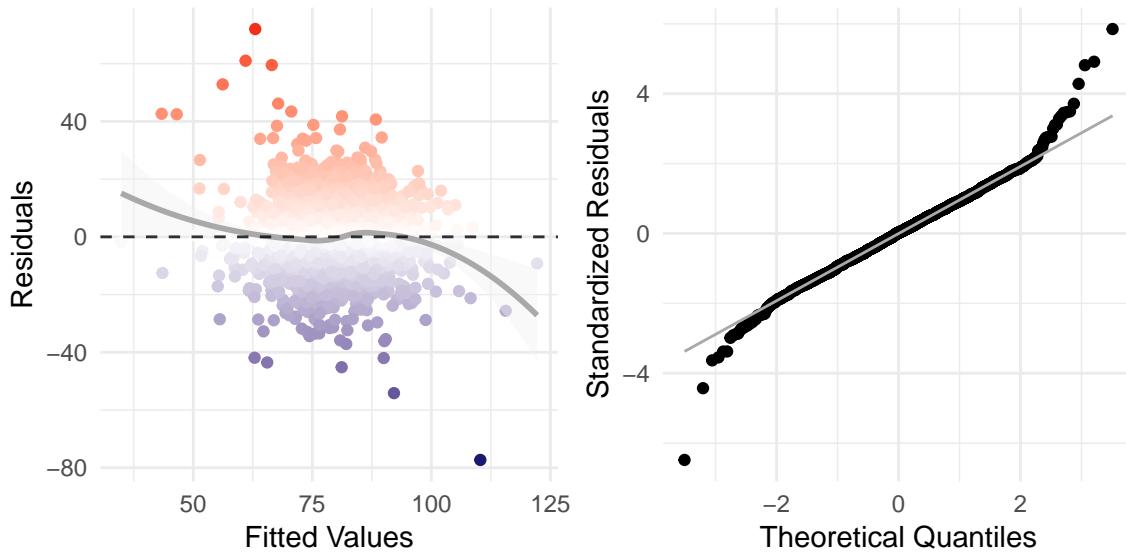


Figure 13: Model 3: Residual Plot and Q-Q Plot

the variance in wins. However, in doing so it treats the batting hits and batting second base runs as drags on wins (with negative coefficients), and pitching hits as buoying wins - which is counterintuitive. While the other coefficients make more intuitive sense, these signs call into question how effectively we can use this model to understand the relationships between the independent variables and wins.

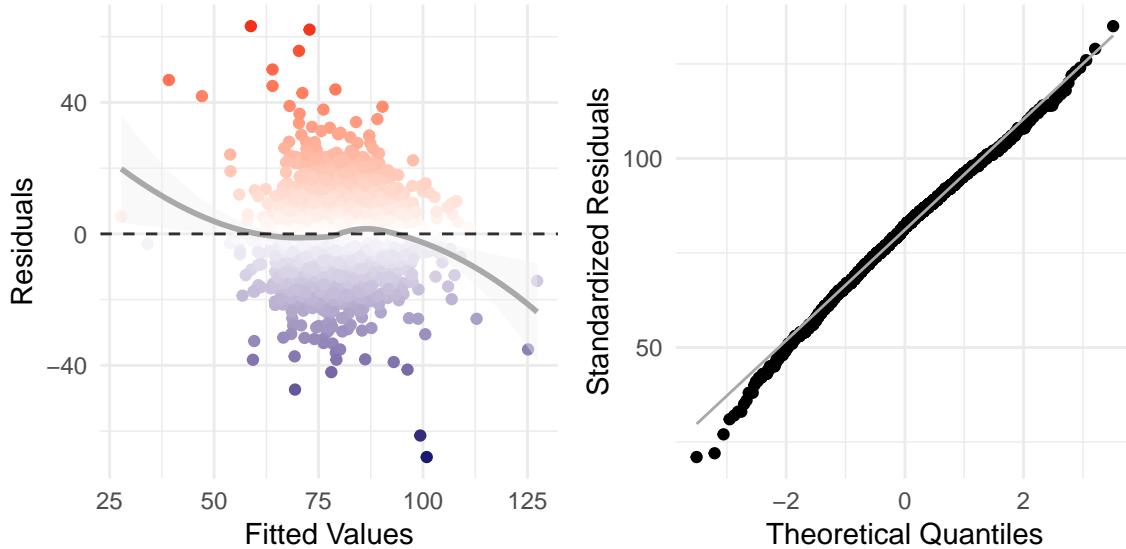


Figure 14: Model 2: Residual Plot and Q-Q Plot

## 4 SELECT MODELS

### 4.1 Comparison of models

In order to determine which model was best suited for determining the number of wins in a season of the three developed, two factors were considered primarily. First and foremost was the  $R^2$  of the models - Model 2 had the highest with an  $R^2$  of 0.3247, while Model 1 had an  $R^2$  of 0.2996 and Model 3 had an  $R^2$  of 0.2956. Secondarily, the coefficients were considered as to whether or not they made sense - the first and second model's both did, whereas the third's did not due to having negative values. It was these two factors that led to the second model ultimately being chosen.

### 4.2 Multi-collinearity

An examination of the partial correlation coefficients between all the independent variables shows that the expected relationships between the offense and defense variables - i.e. batting and pitching homeruns - are there, meet p-value thresholds, and are strong. This finding suggests keeping only one of each pair of variables in a given model; yet, when we tried this, less of the variance in wins was explained by the model. As a result we chose to transform each set of two collinear variables by multiplying them together into one combined variable in our model. This resulted in the highest  $R^2$  value of any of the models we created.

### 4.3 Check Conditions for Least Squares Regression

Our final model's residuals are normally distributed and random. Furthermore, the residuals present with constant variability and no indication of homoscedasticity.

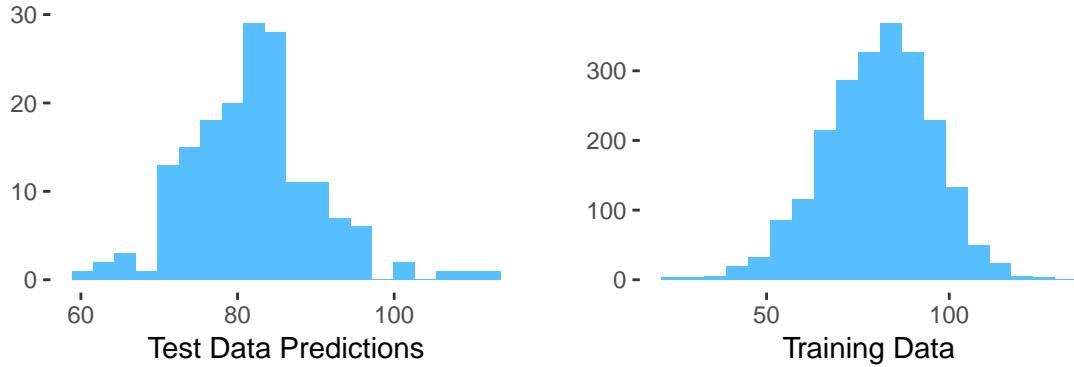


Figure 15: Predictions vs. training data

#### 4.4 F-Statistic

The F-Statistic for the final model is 70.97 on 15 and 2,214 Degrees of Freedom with a p-value less than  $2.2e^{-16}$ . The large F-Statistic indicates we can reject the null hypothesis that the various variables used for the linear model are in no way related to the number of wins a team has in a season. This further supports the model's validity.

#### 4.5 Predictions

We ran predictions on our final model and plotted the distribution next to the distribution from our target in the training data set to compare...

---

## 5 Appendix

The appendix is available as script.R file in Project1 folder. [https://github.com/betsyrosalen/DATA\\_621\\_Business\\_Analyt\\_and\\_Data\\_Mining](https://github.com/betsyrosalen/DATA_621_Business_Analyt_and_Data_Mining)

```
# load data
train <- read.csv('https://raw.githubusercontent.com/silverrainb/data621proj1/master/moneyball-training'
                  stringsAsFactors = F, header = T)
test <- read.csv('https://raw.githubusercontent.com/silverrainb/data621proj1/master/moneyball-evaluation'
                  stringsAsFactors = F, header = T)
# check data
str(train)
str(test)

# remove index
train$INDEX <- NULL
test$INDEX <- NULL

# clean the variable names so it is easier to use
cleanVar <- function(data) {
  name.list <- names(data)
  name.list <- gsub("TEAM_", "", name.list)
  names(data) <- name.list
  data
}

# apply the function
train <- cleanVar(train)
test <- cleanVar(test)

# check data once again
str(train)
str(test)

# Tables and Figures

# Summary Statistics
tbl1 <- describe(train)[,c(2,8,3,5,9,4)]

# Histogram
Hist <- train %>%
  gather() %>%
  ggplot(aes(value)) +
  facet_wrap(~ key, scales = "free") +
  geom_histogram(fill = "#58BFFF") +
  xlab("") +
  ylab("") +
  theme(panel.background = element_blank())
```

```

# Boxplot
melt.train <- melt(train)

outlier.boxplot <- ggplot(melt.train, aes(variable, value)) +
  geom_boxplot(width=.5, fill="#58BFFF", outlier.colour="#58BFFF", outlier.size = 1) +
  scale_y_log10() +
  stat_summary(aes(colour="mean"), fun.y=mean, geom="point",
               size=2, show.legend=TRUE) +
  stat_summary(aes(colour="median"), fun.y=median, geom="point",
               size=2, show.legend=TRUE) +
  coord_flip(ylim = c(0, 2200), expand = TRUE) +
  scale_y_continuous(labels = scales::comma,
                      breaks = seq(0, 2200, by = 200)) +
  labs(colour="Statistics", x="", y="log transformed freq.") +
  scale_colour_manual(values=c("red", "blue")) +
  theme(panel.background=element_blank(), legend.position="top")

# Linearity
linearity <- train %>%
  gather(-TARGET_WINS, key = "var", value = "value") %>%
  ggplot(aes(x = value, y = TARGET_WINS)) +
  geom_point(alpha=0.1) +
  stat_smooth() +
  facet_wrap(~ var, scales = "free", ncol=3) +
  ylab("TARGET_WINS") +
  xlab("") +
  theme(panel.background = element_blank())

# Missing values


```

```

# Remove outliers

max_sd = 5 # change this number to change the threshold for how many standard deviations from the mean to consider an outlier

outliers <- sapply(imputed_train[,-1], function(x) ifelse(x < mean(x)+(sd(x)*max_sd), TRUE, NA))
#outliers <- sapply(imputed_train[,-1], function(x) ifelse(findInterval(x, c(mean(x)-(sd(x)*max_sd),mean(x)+(sd(x)*max_sd))) == 0, TRUE, NA))
imputed_train <- imputed_train[complete.cases(outliers),]

sapply(train.mod, function(x) sum(is.na(x)))

train.mod[, `:=` (BATTING_SO = imputation$BATTING_SO,
                 BASERUN_SB = imputation$BASERUN_SB,
                 BASERUN_CS = imputation$BASERUN_CS,
                 PITCHING_SO = imputation$PITCHING_SO,
                 FIELDING_DP = imputation$FIELDING_DP)]

# Correlations
corr.train <- round(cor(imputed_train),3)
corr.plot <- ggcorrplot::ggcorrplot(corr.train,
                                      type = 'lower',
                                      lab=T,
                                      lab_size=2)

# Feature Engineering
imputed_train$BP_H <- imputed_train$BATTING_H - imputed_train$PITCHING_H
imputed_train$BP_HR <- imputed_train$BATTING_HR - imputed_train$PITCHING_HR
imputed_train$BP_BB <- imputed_train$BATTING_BB - imputed_train$PITCHING_BB
imputed_train$BP_SO <- imputed_train$BATTING_SO - imputed_train$PITCHING_SO

# Model 1
model_exp <- lm(TARGET_WINS ~ BATTING_H + BATTING_HR + BATTING_SO + FIELDING_E +
                  PITCHING_SO + BASERUN_SB + PITCHING_HR + BATTING_BB + BATTING_2B +
                  BATTING_3B + FIELDING_DP + PITCHING_BB + PITCHING_H ,
                  data = imputed_train)

model_exp2 <- lm(TARGET_WINS ~ BATTING_HR + BATTING_SO + FIELDING_E +
                  PITCHING_SO + BASERUN_SB + PITCHING_HR + BATTING_BB +
                  BATTING_3B + FIELDING_DP + PITCHING_BB + PITCHING_H ,
                  data = imputed_train)

figure6 <- stripchart(data.frame(scale(imputed_train)), method ="jitter", las=2,
                       vertical=TRUE)

model_exp3 <- lm(TARGET_WINS ~ BATTING_HR + BATTING_SO + FIELDING_E +
                  PITCHING_SO + BASERUN_SB + PITCHING_HR + BATTING_BB +
                  BATTING_3B + FIELDING_DP + PITCHING_BB + log(PITCHING_H) ,
                  data = imputed_train)

residplot_exp <- ggplot(data = model_exp3,
                         aes(x = .fitted,
                             y = .resid)) +
  geom_point(aes(y = .resid,
                 color = .resid)) +

```

```

        scale_color_gradient2(low = "midnightblue",
                               mid = 'white',
                               high = 'red2') +
      stat_smooth(method = 'loess',
                  se = TRUE,
                  fill = 'gray95',
                  color = 'darkgray') +
      geom_hline(yintercept = 0,
                  col = "black",
                  linetype = "dashed",
                  alpha = .8,
                  size = .5) +
      guides(color = FALSE) +
      labs(x = 'Fitted Values',
            y = 'Residuals') +
      theme_minimal() +
      scale_y_continuous(labels = scales::comma) +
      scale_x_continuous(labels = scales::comma) +
      theme(plot.title = element_text(hjust = .5))

qqplot_exp <- ggplot(data = imputed_train, aes(sample = TARGET_WINS)) +
  stat_qq(size = 1.5) +
  stat_qq_line(color = 'darkgray') +
  labs(x = "Theoretical Quantiles",
       y = "Standardized Residuals") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = .5))

pred_exp <- predict(model_exp3, test)

# Model 2

plot1 <- ggplot(imputed_train, aes(x = BATTING_HR, y = PITCHING_HR)) +
  geom_point(alpha=0.1) +
  stat_smooth() +
  theme(panel.background = element_blank())

plot2 <- ggplot(imputed_train, aes(x = BATTING_HR, y = BATTING_SO)) +
  geom_point(alpha=0.1) +
  stat_smooth() +
  theme(panel.background = element_blank())

plot3 <- ggplot(imputed_train, aes(x = BATTING_BB, y = PITCHING_BB)) +
  geom_point(alpha=0.1) +
  stat_smooth() +
  theme(panel.background = element_blank())

plot4 <- ggplot(imputed_train, aes(x = BATTING_SO, y = PITCHING_SO)) +
  geom_point(alpha=0.1) +
  stat_smooth() +
  theme(panel.background = element_blank())

plot5 <- ggplot(imputed_train, aes(x = BATTING_SO, y = PITCHING_HR)) +
  geom_point(alpha=0.1) +

```

```

stat_smooth() +
  theme(panel.background = element_blank())

plot6 <- ggplot(imputed_train, aes(x = BASERUN_SB, y = BASERUN_CS)) +
  geom_point(alpha=0.1) +
  stat_smooth() +
  theme(panel.background = element_blank())

lm_data <- imputed_train[,-c(9,11:13,16:19)]
lm_data <- as.data.frame(lm_data)

fig8 <- lm_data %>%
  gather(~TARGET_WINS, key = "var", value = "value") %>%
  ggplot(aes(x = value, y = TARGET_WINS)) +
  geom_point(alpha=0.1) +
  stat_smooth() +
  facet_wrap(~ var, scales = "free", ncol=3) +
  xlab("") +
  ylab("TARGET_WINS") +
  theme(panel.background = element_blank())

Histograms <- lm_data %>%
  gather() %>%
  ggplot(aes(value)) +
  facet_wrap(~ key, scales = "free") +
  geom_histogram(fill = "#58BFFF") +
  xlab("") +
  ylab("") +
  #ggtitle("Histograms") +
  theme(panel.background = element_blank())

# Log Transform Data
to_log <- c("BASERUN_SB", "BATTING_3B", "FIELDING_E", "PITCHING_H")
log_lm_data <- lm_data
log_lm_data[,to_log] <- log(log_lm_data[,to_log])

fig10 <- log_lm_data[,c(to_log, "TARGET_WINS")] %>%
  gather(~TARGET_WINS, key = "var", value = "value") %>%
  ggplot(aes(x = value, y = TARGET_WINS)) +
  geom_point(alpha=0.1) +
  stat_smooth() +
  facet_wrap(~ var, scales = "free", ncol=2) +
  xlab("") +
  ylab("TARGET_WINS") +
  theme(panel.background = element_blank())

fig11 <- log_lm_data[,to_log] %>%
  gather() %>%
  ggplot(aes(value)) +
  facet_wrap(~ key, scales = "free", ncol=4) +
  geom_histogram(fill = "#58BFFF") +
  xlab("") +
  ylab("") +
  #ggtitle("Histograms") +

```

```

theme(panel.background = element_blank())

# Model 2 - first model
# Basic linear model with all variables
lm1 <- lm(TARGET_WINS ~ BATTING_H + BATTING_2B + log(BATTING_3B) + BATTING_HR + BATTING_BB + BATTING_SO)
lm_summary <- summary(lm1)

# All Subsets Regression from leaps package
leaps <- regsubsets(x=log_lm_data[,-1], y=log_lm_data[,1], nbest=3)
# plot a table of models showing variables in each model.
# models are ordered by the selection statistic.
leaps_plot <- plot(leaps, scale="r2")

# Scale all the predictor variables
z_train <- data.frame(cbind(lm_data[,1], sapply(lm_data[,-1], scale)))
# Linear model using all scaled predictors
colnames(z_train)[1] <- "TARGET_WINS"
scaled_lm <- lm(TARGET_WINS ~ BATTING_H + BATTING_2B + BATTING_3B + BATTING_HR + BATTING_BB + BATTING_SO)
scaled_lm_summary <- summary(scaled_lm)

#nullmod
nullmod <- lm(TARGET_WINS ~ 1, lm_data)
anova(nullmod, lm1)
### Test one predictor
lm2 <- lm(TARGET_WINS ~ ., lm_data[, -2])
anova(lm2, lm1)
### Test one predictor
lm3 <- lm(TARGET_WINS ~ ., lm_data[, -3])
anova(lm3, lm1)

#Testing a subspace

all_data <- imputed_train[,-c(16:19)]
lm4 <- lm(TARGET_WINS ~ I(BATTING_HR+PITCHING_HR)+I(BATTING_BB+PITCHING_BB)+I(BATTING_SO+PITCHING_SO)+BATTING_H+BATTING_2B+log(BATTING_3B)+log(BASERUN_SB)+BASERUN_CS+log(PITCHING_H)+log(FIELDING_E)+FIELDING_DP, all_data)
lm4_summary <- summary(lm4)

mod_1 <- lm(TARGET_WINS ~ ., imputed_train)
step <- stepAIC(mod_1, direction="both")
#step$anova # display results

#changed to BATTING_BB*PITCHING_BB, BATTING_SO*PITCHING_SO, BATTING_H*log(PITCHING_H)
mod_3 <- lm(TARGET_WINS ~ BATTING_3B + BATTING_HR + BATTING_BB*PITCHING_BB + BATTING_SO*PITCHING_SO + BASERUN_SB + BASERUN_CS + BATTING_H*log(PITCHING_H) + log(FIELDING_E) + FIELDING_DP, imputed_train)

mod_3_summary <- summary(mod_3)

# Model 3

Histograms.mod3 <- train.mod[,c(10, 12:14)] %>%
  gather() %>%
  ggplot(aes(value)) +

```

```

facet_wrap(~ key, scales = "free", nrow = 1) +
geom_histogram(fill = "#58BFFF") +
theme(panel.background = element_blank())

logtransform_lm <- lm(TARGET_WINS ~
  BATTING_H
  + BATTING_2B
  + BATTING_3B
  + BATTING_HR
  + BATTING_BB
  + BATTING_SO
  + BASERUN_SB
  #+ BASERUN_CS
  + log(PITCHING_H)
  #+ log(PITCHING_HR + .0001) # p-value around .16 as log .27 w/o so remove
  + PITCHING_BB
  #+ log(PITCHING_SO + .0001) # p-value around .5 whether or not log transform
  + FIELDING_E
  + FIELDING_DP
, data = imputed_train)

logtransform_lm_summary <- summary(logtransform_lm)

#logtransform_lm_summary

# Model 3: append predictions and residuals
imputed_train$logtransform_pred <- predict(logtransform_lm)
imputed_train$logtransform_resid <- residuals(logtransform_lm)

# Model 3: residual plot
logtransform_residplot <- ggplot(data = logtransform_lm, # for each model update chart object name and c
  aes(x = .fitted,
      y = .resid)) +
  geom_point(aes(y = .resid,
                 color = .resid)) +
  scale_color_gradient2(low = "midnightblue",
                        mid = 'white',
                        high = 'red2') +
  stat_smooth(method = 'loess',
              se = TRUE,
              fill = 'gray95',
              color = 'darkgray') +
  geom_hline(yintercept = 0,
             col = "black",
             linetype = "dashed",
             alpha = .8,
             size = .5) +
  guides(color = FALSE) +
  labs(x = 'Fitted Values',
       y = 'Residuals') +
  theme_minimal() +
  scale_y_continuous(labels = scales::comma) +
  scale_x_continuous(labels = scales::comma) +

```

```

theme(plot.title = element_text(hjust = .5))

# Model 3: QQ-plot residuals
logtransform_qqplot <- ggplot(logtransform_lm, aes(sample = .stdresid)) + # for each model update char
  stat_qq(size = 1.5) +
  stat_qq_line(color = 'darkgray') +
  labs(x = "Theoretical Quantiles",
       y = "Standardized Residuals") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = .5))

#Multi-colinearity
collintest_df <- imputed_train[,3:15]
# subset dataframe for independent variables
collintest_stats <- pcor(collintest_df, method = 'pearson')

# Check Conditions for Least Squares Regression
residplot_exp <- ggplot(data = mod_3,
                         aes(x = .fitted,
                             y = .resid)) +
  geom_point(aes(y = .resid,
                 color = .resid)) +
  scale_color_gradient2(low = "midnightblue",
                        mid = 'white',
                        high = 'red2') +
  stat_smooth(method = 'loess',
              se = TRUE,
              fill = 'gray95',
              color = 'darkgray') +
  geom_hline(yintercept = 0,
             col = "black",
             linetype = "dashed",
             alpha = .8,
             size = .5) +
  guides(color = FALSE) +
  labs(x = 'Fitted Values',
       y = 'Residuals') +
  theme_minimal() +
  scale_y_continuous(labels = scales::comma) +
  scale_x_continuous(labels = scales::comma) +
  theme(plot.title = element_text(hjust = .5))

qqplot_exp <- ggplot(data = lm_data, aes(sample = TARGET_WINS)) +
  stat_qq(size = 1.5) +
  stat_qq_line(color = 'darkgray') +
  labs(x = "Theoretical Quantiles",
       y = "Standardized Residuals") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = .5))

# Predictions
predictions <- round(predict(mod_3, test))

p1.pred <- ggplot(data.frame(predictions), aes(predictions)) +

```

```
geom_histogram(fill = "#58BFFF", bins = 20) +
  xlab("Test Data Predictions") +
  ylab("") +
  theme(panel.background = element_blank())

p2.pred <- ggplot(lm_data, aes(TARGET_WINS)) +
  geom_histogram(fill = "#58BFFF", bins = 20) +
  xlab("Training Data") +
  ylab("") +
  theme(panel.background = element_blank())
```